The background features a light beige color with several large, semi-transparent hexagonal shapes. Inside these hexagons, there are abstract geometric patterns: one shows white cubes connected by thin lines, and another shows a network of white dots connected by lines. The overall aesthetic is clean and modern.

Mateusz Wasiluk  
ROBIN Lab Meeting  
09.02.2023

# Bayesian Neural Networks for Continual Learning

# Plan for today

- The principles of Bayesian learning
- Inference methods for BNNs
- Mean field approximation in practice
- Introduction to normalizing flows

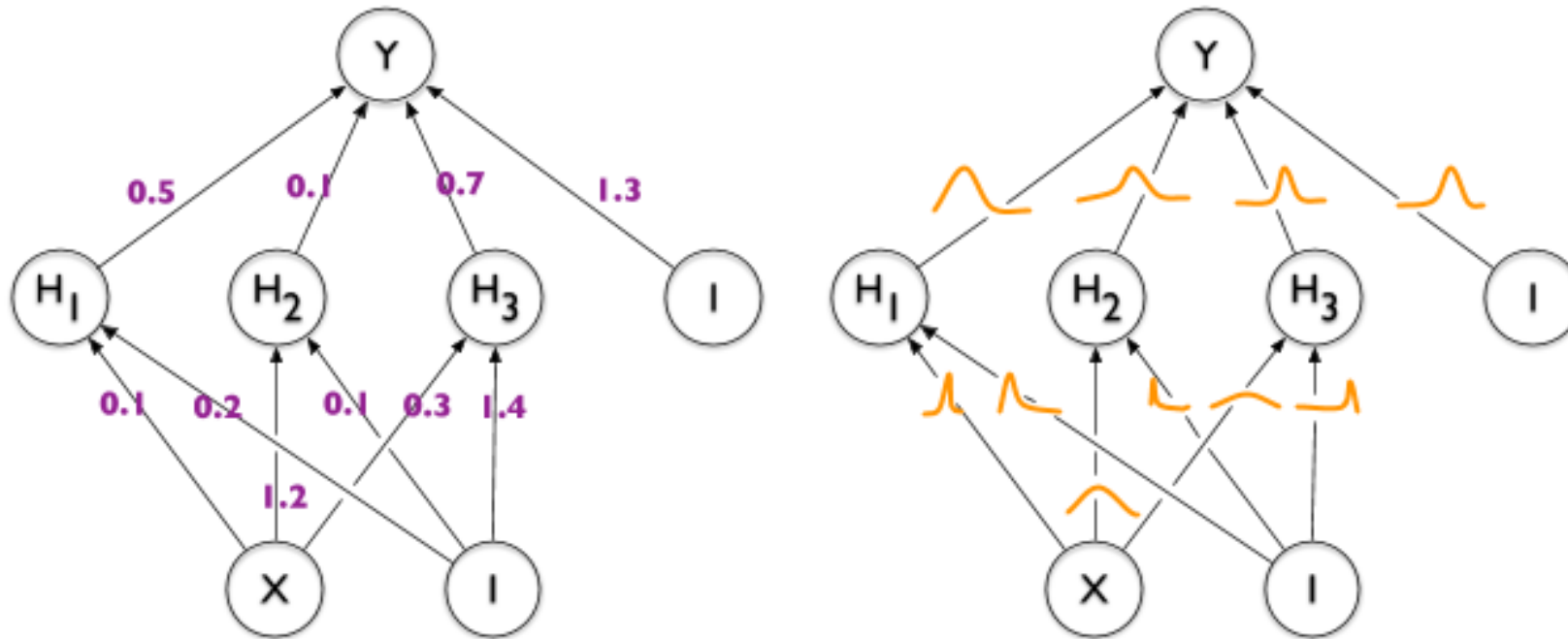
# New observations influence our model of the world

$$\textit{Posterior} \propto \textit{prior} * \textit{likelihood}$$

The new model depends on the old beliefs and the new observations *in the light of the old beliefs*

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)} = \frac{P(D, H)}{\int_H P(D, H')dH'}$$

# A BNN is any stochastic neural network trained with Bayesian Inference



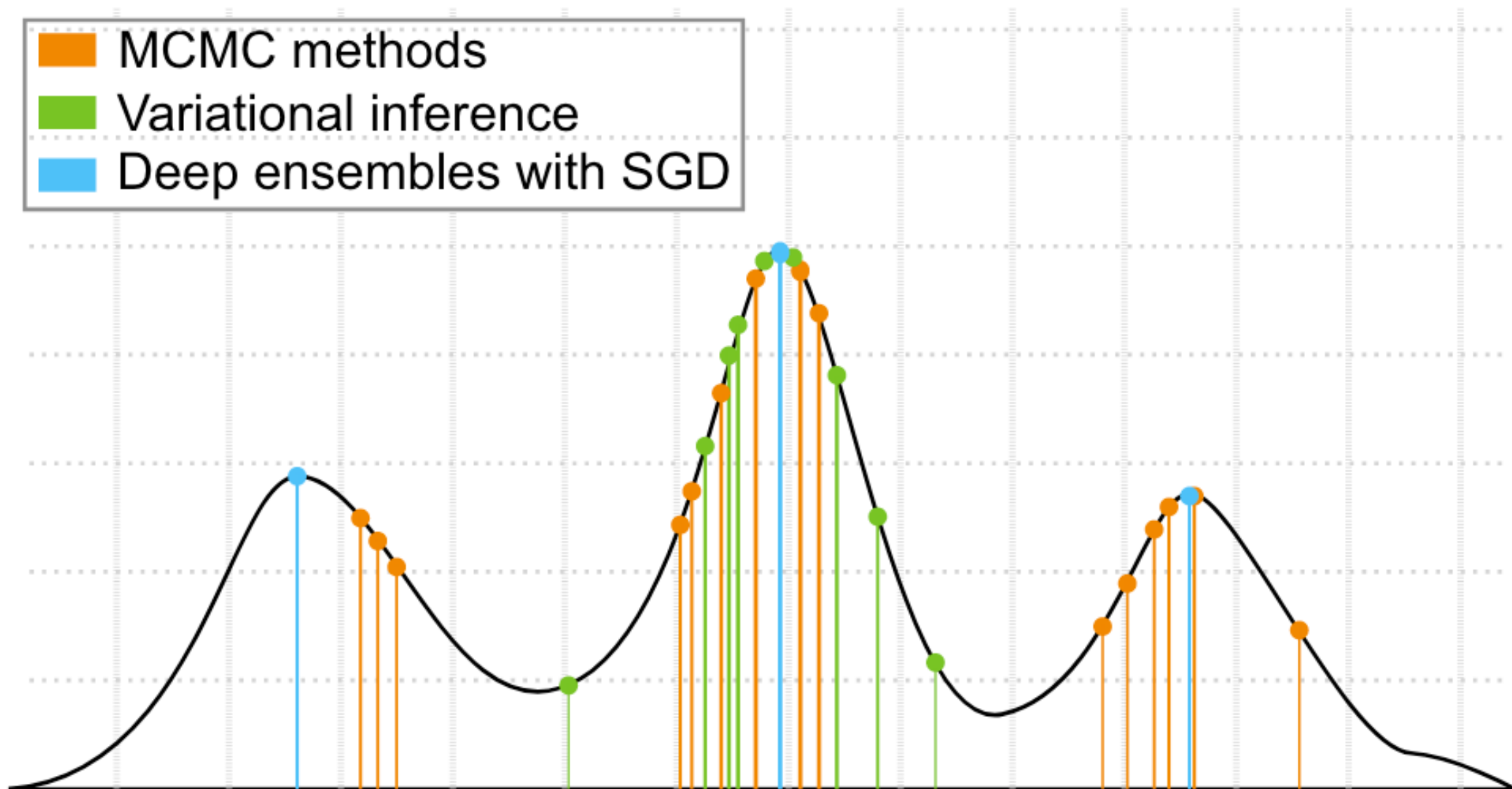
# Most popular methods of inference are MCMC and VI

## Markov Chain Monte Carlo

- Directly samples the posterior
- Can explore different modes
- Does not scale well

## Variational Inference

- Uses a parametric approximation of the posterior
- Aims at minimizing the KL-divergence
- Scales for large models with a limited expressive power

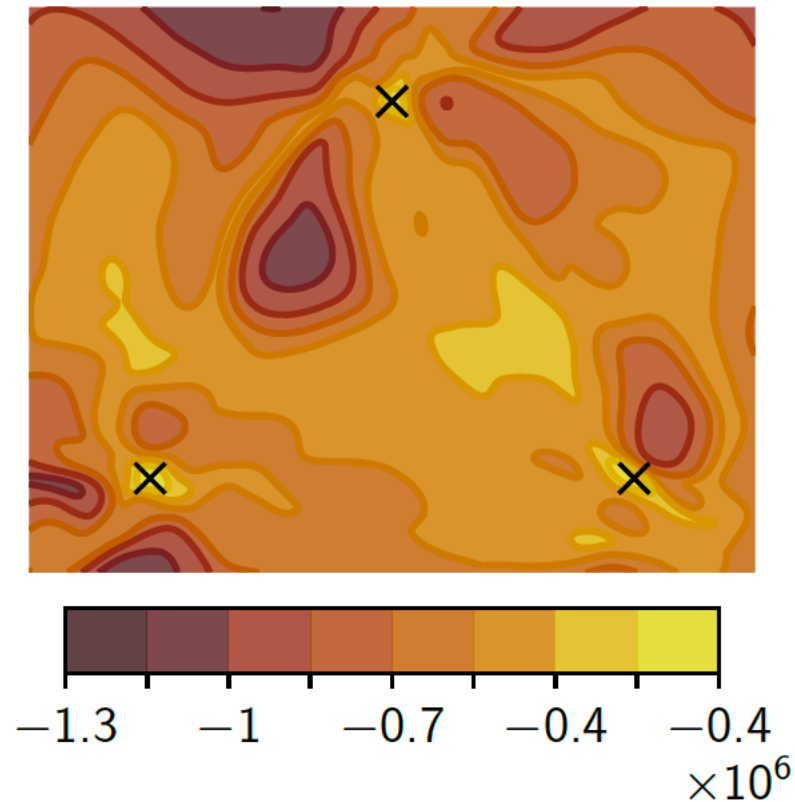


# Training a variational network requires maximizing the ELBO

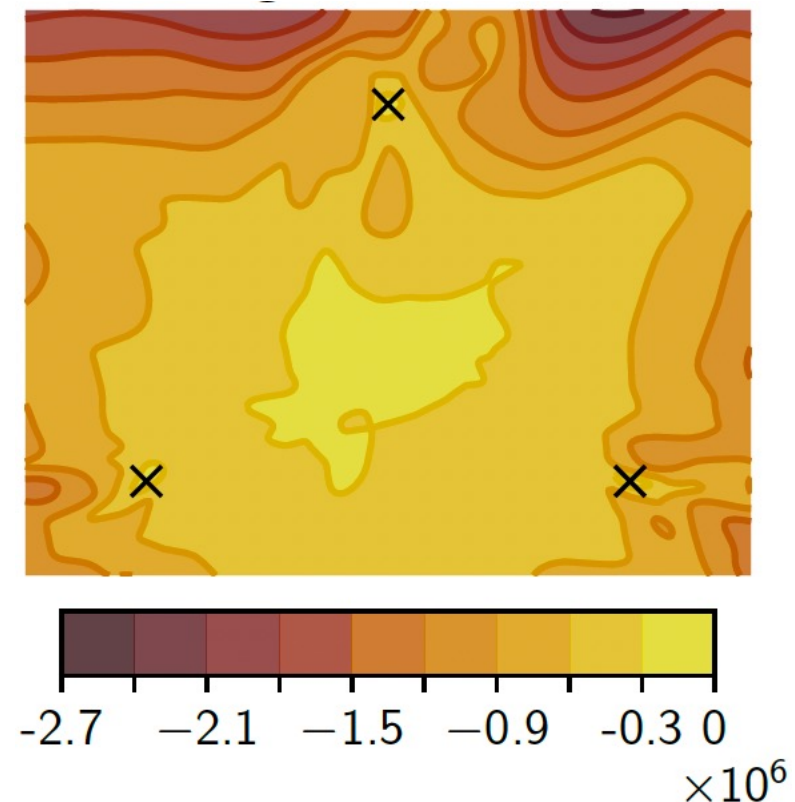
$$\underbrace{\underbrace{\log P(D)}_{\text{evidence}} - \underbrace{D_{KL}[Q(\theta)||P(\theta|D)]}_{\text{variational-posterior divergence}}}_{\text{ELBO (training loss)}} = \underbrace{E_{\theta \sim Q}[\log P(Y|X, \theta)]}_{\text{classification loss}} - \underbrace{D_{KL}[Q(\theta)|P(\theta)]}_{\text{variational-prior divergence}}$$

# Mean-field assumption is too strong for complex datasets

Log-posterior, CIFAR-10, samples from one HMC chain

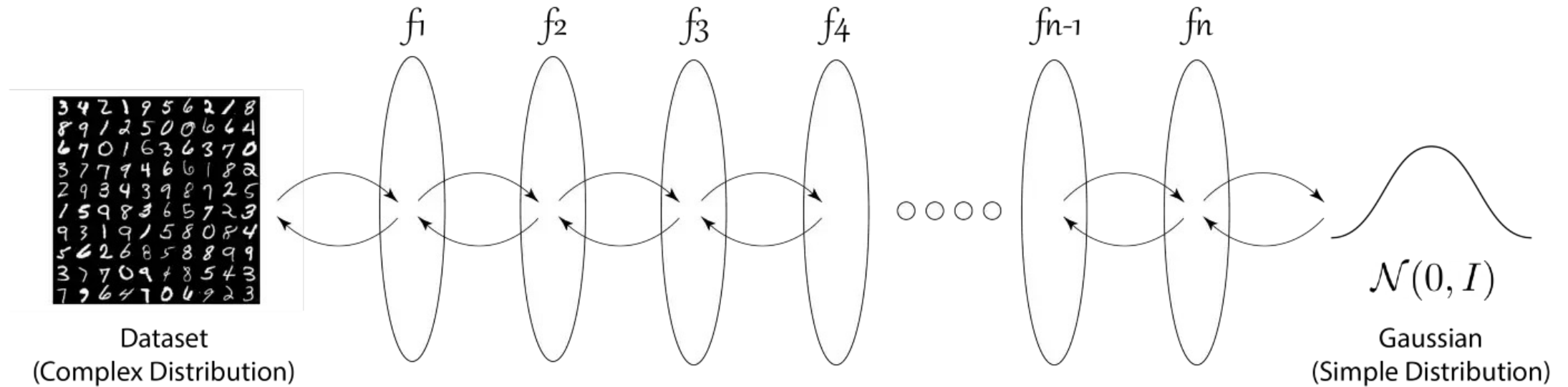


Log-posterior, CIFAR-10, samples from independent HMC chains

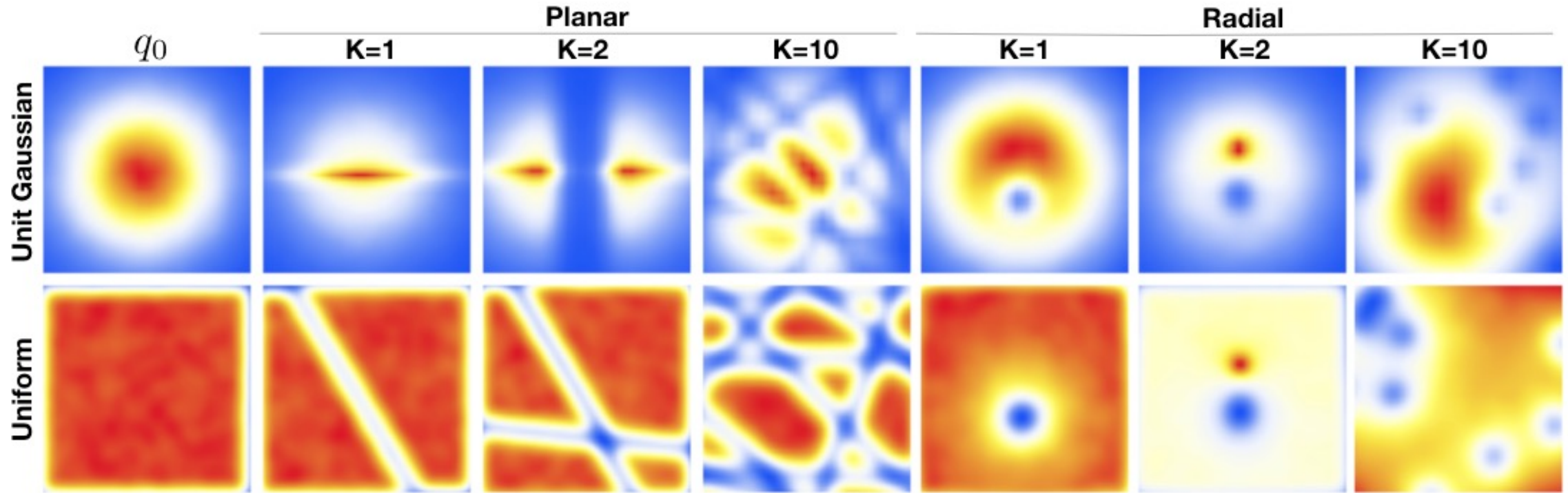




# Normalizing flows transform simple distributions to more complex ones



# Example: "geometrical" flows

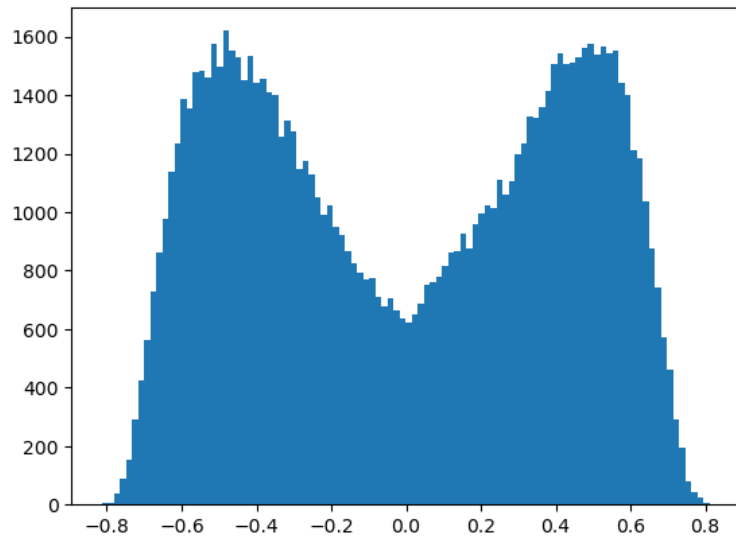


# Want to understand more?

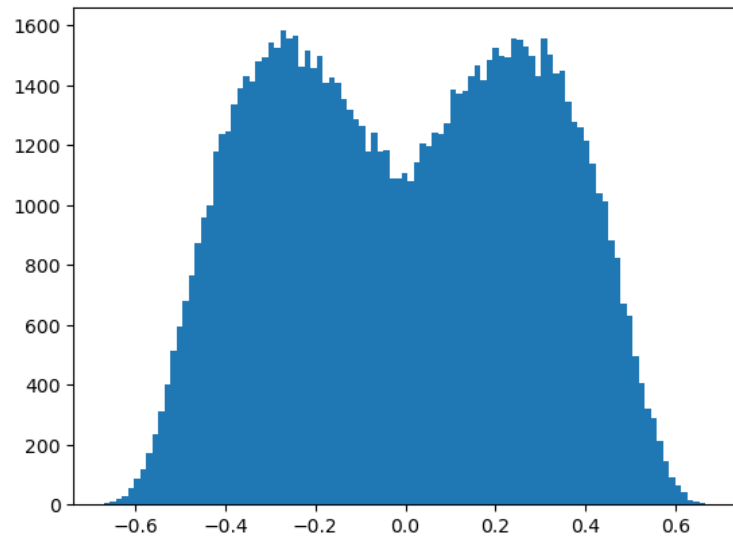
- Rezende, D., Mohamed, S. (2015). Variational Inference with Normalizing Flows
- Jospin, L.V., Laga, H., Boussaid, F., Buntine, W. and Bennamoun, M., 2022. Hands-on Bayesian neural networks—A tutorial for deep learning users.
- Izmailov, P., Vikram, S., Hoffman, M.D., Wilson, A.G.G.. (2021). What Are Bayesian Neural Network Posteriors Really Like?

# Backup

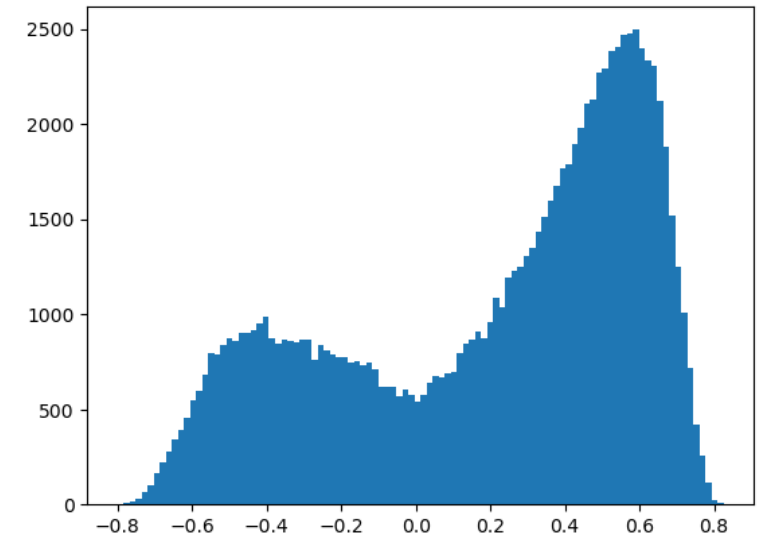
# Example: softsign flow's effect on the Gaussian



$N(0, 1)$

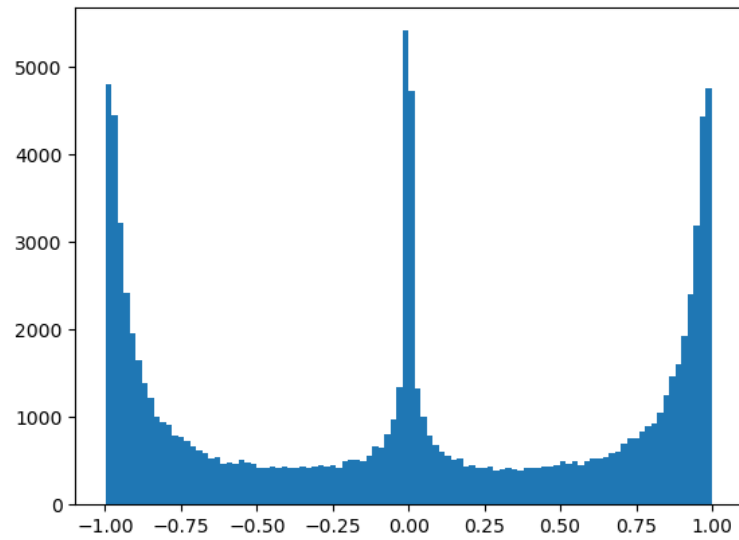


$N(0, 0.5)$

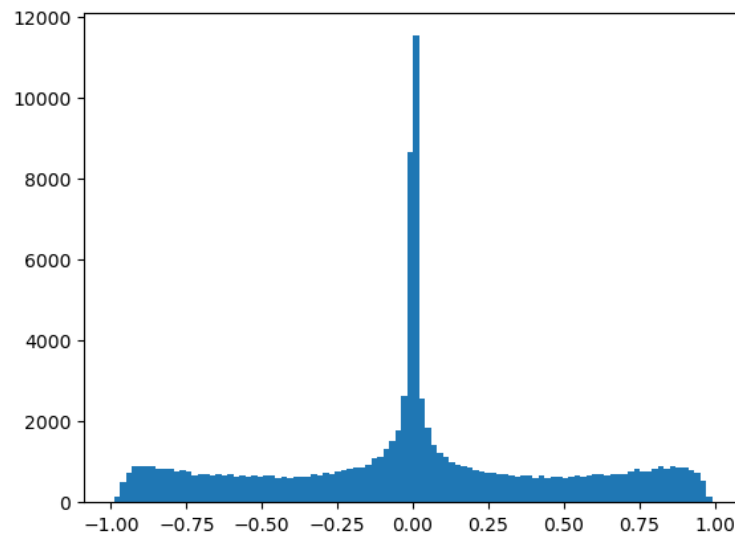


$N(0.5, 1)$

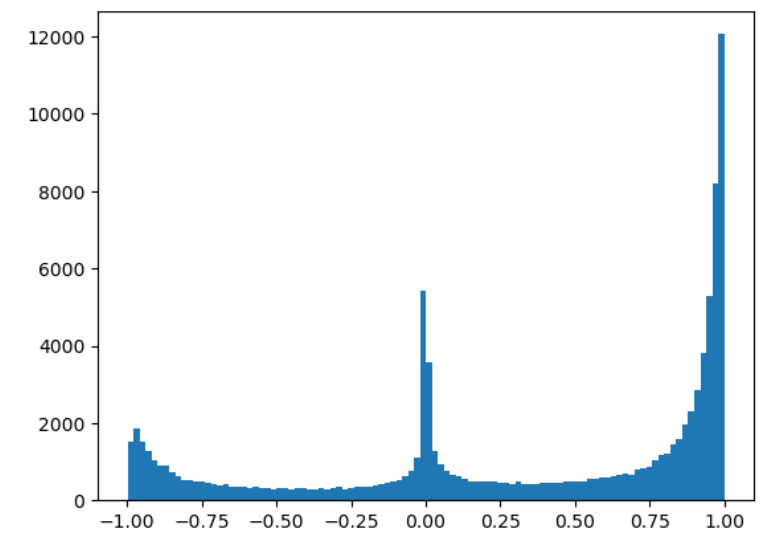
# Example: softsign->scale(5)->power(3)



$N(0, 1)$



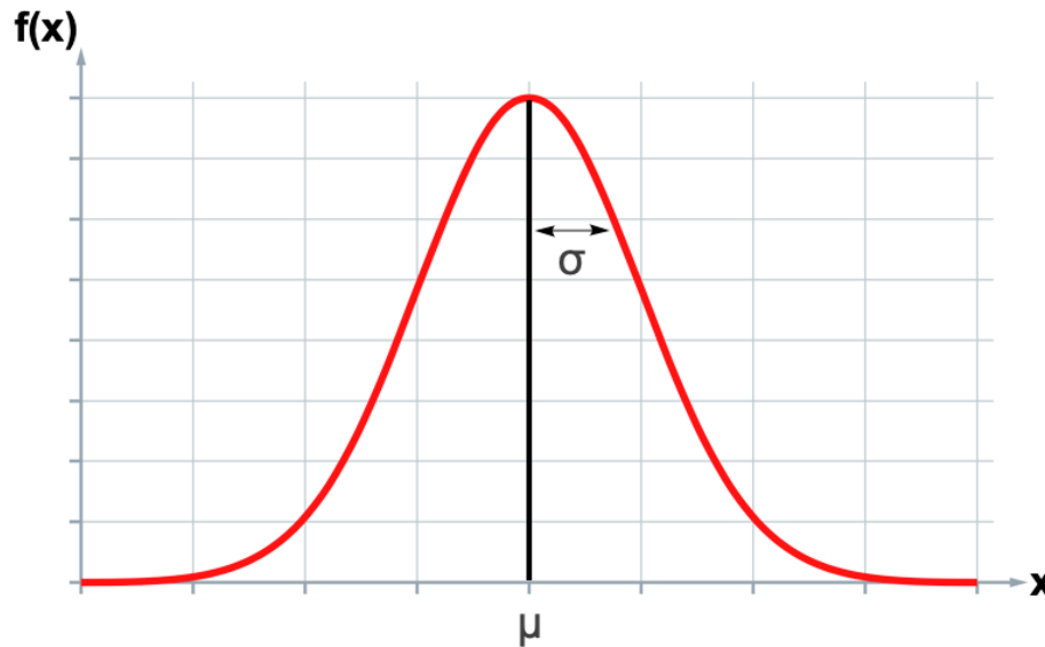
$N(0, 0.5)$



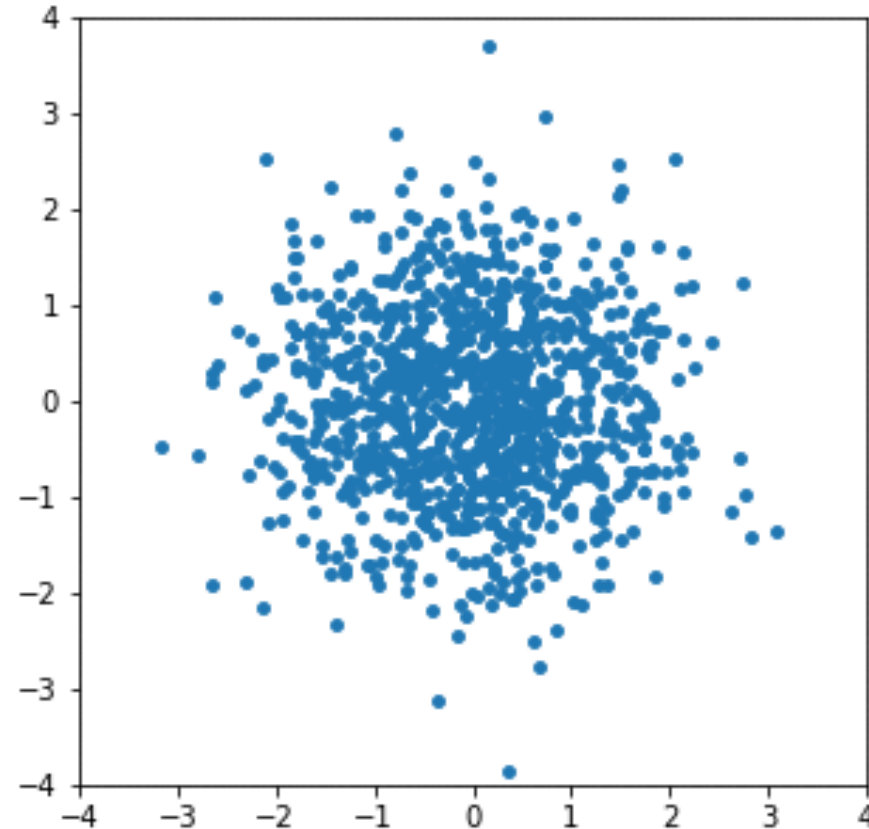
$N(0.5, 1)$

The “reparametrization trick” allows SGD methods to pass through stochastic nodes

$$w = \mu + \sigma * \epsilon, \quad \epsilon \sim N(0, 1)$$



Visualization how samples from a Gaussian base distribution are transformed through 8 layers of Real NVP to match the target distribution:





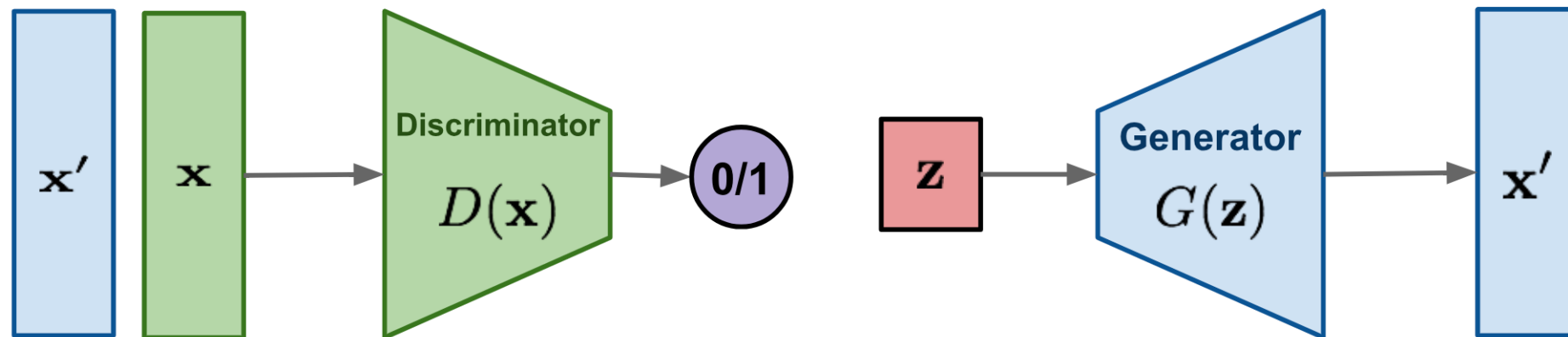
# Posterior formula for a BNN – choosing the network's parameters

$$p(\boldsymbol{\theta}|D) = \frac{p(D_{\mathbf{y}}|D_{\mathbf{x}}, \boldsymbol{\theta})p(\boldsymbol{\theta})}{\int_{\boldsymbol{\theta}} p(D_{\mathbf{y}}|D_{\mathbf{x}}, \boldsymbol{\theta}')p(\boldsymbol{\theta}')d\boldsymbol{\theta}'} \propto p(D_{\mathbf{y}}|D_{\mathbf{x}}, \boldsymbol{\theta})p(\boldsymbol{\theta})$$

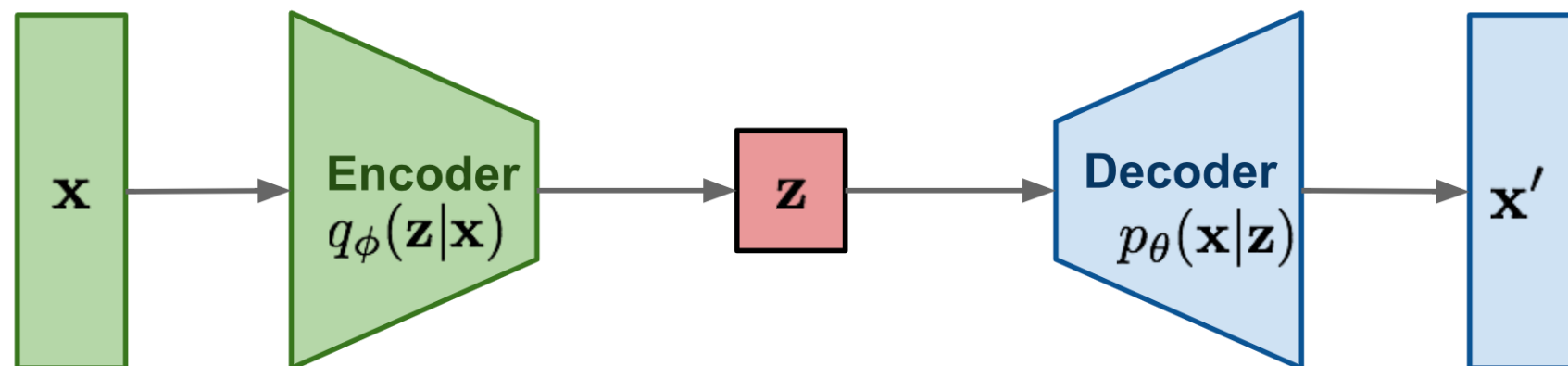
# Predicting with a BNN – the marginal

$$p(\mathbf{y}|\mathbf{x}, D) = \int_{\boldsymbol{\theta}} p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}') p(\boldsymbol{\theta}'|D) d\boldsymbol{\theta}'.$$

**GAN:** minimax the classification error loss.



**VAE:** maximize ELBO.



**Flow-based generative models:** minimize the negative log-likelihood

